



PATENT ABSTRACTS OF JAPAN

(11) Publication number: **2002189492 A**(43) Date of publication of application: **05.07.02**

(51) Int. Cl.

G10L 15/06**G10L 15/10****G10L 15/14****G10L 15/20****G10L 21/02**(21) Application number: **2000385212**(22) Date of filing: **19.12.00**(71) Applicant: **SHARP CORP**(72) Inventor: **YAMAGUCHI KOICHI
HACHIMAN YOICHIRO**

(54) **APPARATUS AND METHOD FOR EXTRACTING
SPEAKER'S FEATURE, SPEECH RECOGNITION
DEVICE, SPEECH SYNTHESIZER, AND
PROGRAM RECORDING MEDIUM**

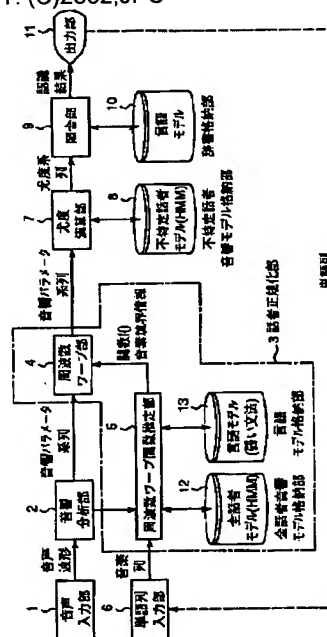
depending on the contents of utterance data.

COPYRIGHT: (C)2002,JPO

(57) Abstract:

PROBLEM TO BE SOLVED: To stably extract a speaker's feature without depending on the contents of utterance data.

SOLUTION: A frequency function estimation part 5 is provided with a phonemic limit estimation part, a frequency measuring part, and a mode extracting part. During learning, the frequency function estimation part 5 estimates phonemic limit information, and performs the maximum likelihood estimation of the coefficient α of the frequency warping function f of each sample about a phonemic section selected on the basis of the phonemic limit information. Furthermore, the frequency function estimation part 5 determines a distribution function $H(\alpha)$ which represents a frequency distribution about the coefficient α of each sample, and estimates a coefficient α which provides a mode value as the optimal coefficient of the frequency warping function f . Consequently, a correct frequency warping function f can be estimated even when a plurality of peaks are present in the frequency distribution, and the speaker's feature can stably be extracted without



(11)特許出願公開番号
特開2002-189492
(P2002-189492A)

(43)公開日 平成14年7月5日(2002.7.5)

(51)Int.Cl. ⁷	識別記号	F I	テーマコード(参考)
G 1 0 L 15/06		G 1 0 L 3/00	5 2 1 S 5 D 0 1 5
15/10			5 3 1 J
15/14			5 3 1 C
15/20			5 3 5 B
21/02		3/02	3 0 1 A
審査請求 未請求 請求項の数10 O L (全 19 頁)			

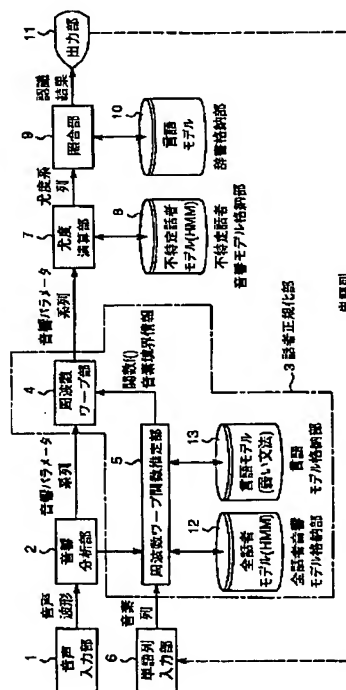
(21)出願番号	特願2000-385212(P2000-385212)	(71)出願人	000005049 シャープ株式会社 大阪府大阪市阿倍野区長池町22番22号
(22)出願日	平成12年12月19日(2000. 12. 19)	(72)発明者	山口 耕市 大阪府大阪市阿倍野区長池町22番22号 シャープ株式会社内
		(72)発明者	八幡 洋一郎 大阪府大阪市阿倍野区長池町22番22号 シャープ株式会社内
		(74)代理人	100062144 弁理士 青山 葆 (外1名) Fターム(参考) 5D015 AA02 BB02 EE03 HH23

(54) 【発明の名称】 話者特徴抽出装置および話者特徴抽出方法、音声認識装置、音声合成装置、並びに、プログラム記録媒体

(57) 【要約】

【課題】 発声データの内容に依存せず安定して話者特徴を抽出する。

【解決手段】 周波数関数推定部5には音素境界推定部、傾度計測部およびモード抽出部を備えている。そして、学習時には、音素境界情報を推定し、この音素境界情報に基づいて選択した音素区間に関して、各サンプル毎に、周波数ワーピング関数 f の係数 α を最尤推定する。さらに、各サンプル毎の係数 α に関する傾度分布を表す分布関数 $H(\alpha)$ を求め、最傾値を与える係数 α を周波数ワーピング関数 f の最適係数として推定する。したがって、上記傾度分布に複数のピークが存在する場合でも正確な周波数ワーピング関数 f を推定でき、発声データの内容に依存せず安定して話者特徴を抽出できる。



【特許請求の範囲】

【請求項 1】 入力話者の音声から、標準話者の音声スペクトルに対して上記入力話者の音声スペクトルの周波数軸を伸縮する際の周波数伸縮関数を話者特徴として抽出する話者特徴抽出装置において、

所定の音声単位毎に、上記標準話者の音響モデルに対して、上記入力話者の音声サンプルの尤度あるいは音響モデルを上記入力話者の音声サンプルに話者適応させた話者適応音響モデルの尤度を最大にするという基準に従って、上記周波数伸縮関数を最尤推定し、この推定された上記周波数伸縮関数の集合の頻度分布を求める頻度計測手段と、

上記頻度分布に基づいて、最大頻度を有する周波数伸縮関数を話者特徴として抽出するモード抽出手段を備えたことを特徴とする話者特徴抽出装置。

【請求項 2】 請求項 1 に記載の話者特徴抽出装置において、

上記入力話者の音声サンプルから音響モデルを用いたビタビアルゴリズムによって音素境界情報を推定する音素境界情報推定手段を備えて、

上記頻度計測手段は、上記音素境界情報に基づいて、有音無音の別および調音点の位置に従って、上記最尤推定を行う際に用いる上記入力話者の音声サンプルの音声区間を限定する機能を有していることを特徴とする話者特徴抽出装置。

【請求項 3】 請求項 1 に記載の話者特徴抽出装置において、

音響モデルの各状態が、有音無音の別および調音点の位置に従って予め設定された上記最尤推定を行う音声区間に属しているか否かを判別し、属している状態に関して、上記音響モデルを入力音声サンプルに話者適応させて話者適応音響モデルを作成する話者適応モデル作成手段を備えて、

上記頻度計測手段は、上記標準話者音響モデルに対して上記話者適応音響モデルを用いて上記最尤推定を行うようになっていることを特徴とする話者特徴抽出装置。

【請求項 4】 請求項 1 乃至請求項 3 の何れか一つに記載の話者特徴抽出装置において、

上記モード抽出手段は、上記最大頻度を有する周波数伸縮関数が複数存在する場合には、上記頻度分布を混合ガウス分布で表現した場合における当該複数の周波数伸縮関数が属している分布の分散が大きい方の周波数伸縮関数をもって話者特徴する機能を有していることを特徴とする話者特徴抽出装置。

【請求項 5】 請求項 4 に記載の話者特徴抽出装置において、

上記モード抽出手段は、上記標準話者の特徴を表す周波数伸縮関数に近い方の周波数伸縮関数をもって話者特徴とする機能を組み合わせて、上記話者特徴を抽出するようになっていることを特徴とする話者特徴抽出装置。

10

20

30

40

50

【請求項 6】 標準話者の音声スペクトルに対して入力話者の音声スペクトルの周波数軸を伸縮する際の周波数伸縮関数を用いて入力話者の音声スペクトルの周波数軸を伸縮することによって上記入力話者の音声を正規化する正規化手段を有する音声認識装置において、
上記正規化手段は、

請求項 1 乃至請求項 5 の何れか一つに記載の話者特徴抽出装置と、

上記話者特徴抽出装置によって抽出された周波数伸縮関数を用いて、上記入力話者の音声スペクトルの周波数軸を伸縮する周波数ワープ手段で構成されていることを特徴とする音声認識装置。

【請求項 7】 入力話者の音声スペクトルに対して標準話者の音声スペクトルの周波数軸を伸縮する際の周波数伸縮関数を用いて音声のスペクトルの周波数軸を伸縮することによって音響モデルを入力話者に話者適応させる話者適応手段を有する音声認識装置において、

上記話者適応手段は、

請求項 1 乃至請求項 5 の何れか一つに記載の話者特徴抽出装置と、

上記話者特徴抽出装置によって抽出された周波数伸縮関数の逆関数を用いて、上記音響モデルの周波数軸を伸縮する周波数ワープ手段で構成されていることを特徴とする音声認識装置。

【請求項 8】 入力話者の音声スペクトルに対して標準話者の音声スペクトルの周波数軸を伸縮する際の周波数伸縮関数を用いて音声のスペクトルの周波数軸を伸縮することによって、標準話者の音声素片を接続して成る合成音声の声質を発話者の声質に変換する声質変換手段を有する音声合成装置において、

上記声質変換手段は、

請求項 1 乃至請求項 5 の何れか一つに記載の話者特徴抽出装置と、

上記話者特徴抽出装置によって抽出された周波数伸縮関数の逆関数を用いて、上記音声素片の周波数軸を伸縮する周波数ワープ手段で構成されていることを特徴とする音声合成装置。

【請求項 9】 入力話者の音声から、標準話者の音声スペクトルに対して上記入力話者の音声スペクトルの周波数軸を伸縮する際の周波数伸縮関数を話者特徴として抽出する話者特徴抽出方法において、

所定の音声単位毎に、上記標準話者の音響モデルに対して、上記入力話者の音声サンプルの尤度あるいは音響モデルを上記入力話者の音声サンプルに話者適応させた話者適応音響モデルの尤度を最大にするという基準に従って、上記周波数伸縮関数を最尤推定し、
この推定された上記周波数伸縮関数の集合の頻度分布を求め、

上記頻度分布に基づいて、最大頻度を有する周波数伸縮関数を話者特徴として抽出することを特徴とする話者特

微抽出方法。

【請求項 10】 コンピュータを、
請求項 1 における上記頻度計測手段およびモード抽出手段として機能させる話者特徴抽出処理プログラムが記録されたことを特徴とするコンピュータ読出し可能なプログラム記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】 この発明は、標準話者の音声スペクトルに対する入力音声スペクトルの周波数軸の線形伸縮関数を話者特徴として抽出する話者特徴抽出装置および話者特徴抽出方法、その抽出方法を用いた音声認識装置、音声合成装置、並びに、話者特徴抽出処理プログラムを記録したプログラム記録媒体に関する。

【0002】

【従来の技術】 従来より、隠れマルコフモデル(Hidden Markov Model : 以下、HMMと言う)を用いた音声認識方法の開発が近年盛んに行われている。このHMMは、大量の音声データから得られる音声の統計的特徴を確率的にモデル化したものであり、このHMMを用いた音声認識方法の詳細は、中川聖一著「確率モデルによる音声認識」(電子情報通信学会)に詳しい。このHMMに基づく話者適応や話者正規化に関する研究が行われている。ところが、通常の話者正規化技術や話者適応技術においては発声データの内容や量に依存するので、少量の発声データからでは安定した性能向上が非常に難しい。そこで、声道長を用いた話者適応や話者正規化の手法が注目されており、特に声道長に基づく話者正規化が盛んに研究されて効果が出ている。声道長は音声のスペクトルの大まかな特徴を表すパラメータであり、声道長の差は話者間の主な変動要因である。また、声道長は従来の話者適応法に比べて 1 個のパラメータあるいは極めて少ないパラメータで音声の特徴を表現できることから、声道長にはより少量の学習データで効率良く正規化できるというメリットがある。

【0003】 ところで、標準話者の音声パターンに対する入力話者の音声サンプルの尤度を最大にするという基準(最尤推定)で、上記音声サンプルにおける周波数軸の線形伸縮係数 α (声道長正規化係数)を求める方法(ML-VTLN法: Maximum Likelihood Vocal Tract Length Normalization)がある。そして、この声道長正規化係数 α を用いて入力話者の音声サンプルの周波数軸を線形伸縮することで話者正規化を行う技術が提案されている(例えば、AT&T Bell Labs. Li Lee, Richard C. Rose, "Speaker Normalization using Efficient Frequency Warping Procedures", pp. 353-356 ICASSP96 (1996))。また、特開平 11-327592 号公報においては、声道を前室と後室との 2 つの室に分け、入力音声のフォルマント周波数を用いて各室に対応した 2 つの周波数軸線形伸縮係数 α を求め、この 2 つの周波数軸線形伸縮係数 α を用

いて話者正規化する技術が開示されている。

【0004】 尚、上記話者適応は標準となる音響モデルを入力話者に対して適応(つまり正規化)させる技術であり、話者正規化とは表裏一体の関係にある。

【0005】 さらに、音声合成における声質変換に関する従来技術として、音声認識の話者適応技術を用いてスペクトルの写像を行なう方法が提案されている。例えば、ベクトル量子化(VQ)コードブックマッピング法をベースとした話者適応技術を用いる方法(特開平 1-97997 号公報)や、VFS (Vector Field Smoothing) 法をベースとした話者適応技術を用いる方法(橋本誠、樋口宣男: "話者選択と移動ベクトル場平滑化を用いた声質変換のためのスペクトル写像", 信学技報, SP95-1, p. p. 1-8, May 1995)等がある。

【0006】

【発明が解決しようとする課題】 しかしながら、上記従来の声道長に基づく話者適応や話者正規化には、以下のような問題がある。すなわち、声道長に基づく話者適応や話者正規化は極めて少ないパラメータ数で音声の特徴を表現できるとは言うものの、声道長の抽出は発声データの内容や量に大きく左右され易い。したがって、必ずしも少ない学習サンプルから安定して声道長を抽出できるとは限らない。その結果、声道長に基づいて話者正規化や話者適応や話者クラスタリングを行うような音声認識装置においては、性能劣化を招くという問題がある。

【0007】 上記声道長正規化係数 α を求める方法としては、上述したように、学習サンプル全体を対象として最尤推定するML-VTLN法が提案されている。この方法においては、発話の仕方や発話内容によってスペクトルの概形が変動するので、学習サンプルによって最適な声道長正規化係数 α が異なってしまう場合が生ずる。つまり、異なる周波数軸伸縮関数で局所的に最適となるため、声道長正規化係数 α の頻度分布に複数のピークが生じてしまうという多峰性の問題が発生し、真の声道長正規化関数が安定して推定できないことになる。

【0008】 また、声道長の伸縮(周波数軸のワーピング)を線形関数やそれに類似した関数で表現しており、一般に全音素区間に対して周波数ワープを学習・作用するようにしている。そのために、声道長正規化係数 α を用いて話者正規化を行う方法においては、声道長の差の影響を受け難い音素や無音部まで学習および正規化してしまうという問題もある。

【0009】 すなわち、上記従来の声道長正規化係数の推定方法では、正確な声道長正規化係数が推定できなかったり、入力音声サンプルを必要以上に変形してしまったりするために、話者正規化に適用した場合には、認識性能の低下を招くことになるのである。

【0010】 さらに、上記特開平 11-327592 号公報の正規化方法においては、2 つのサンプルから直接声道パラメータを求めるようにしているが、声道パラメ

ータを得るために入力音声のフォルマント周波数を用いている。ところが、一般的にフォルマント周波数を全自動で求めることは困難であり、上記特開平11-327592号公報に開示された線形伸縮係数を用いた話者正規化方法では、実時間性に欠けるという問題がある。

【0011】また、上記話者適応においては少ない発声データから音響モデルを精度良く適応できないために、誤り率を半減させるためには数十単語以上の発声データが必要となり、学習話者に負担を強いることになるという問題がある。さらに、音響モデルの適応に声道長伸縮関数を用いる場合には、上述した話者正規化の場合と同様の問題が発生する。また、音声合成における声質変換の場合にも、同様に少ない発声データからは精度良く声質が得られないという問題がある。

【0012】そこで、この発明の目的は、少ない発声データから発声データの内容に依存せずに安定して話者特徴を抽出できる話者特徴抽出装置および話者特徴抽出方法、その抽出方法を用いた音声認識装置、音声合成装置、並びに、話者特徴抽出処理プログラムを記録したプログラム記録媒体を提供することにある。

【0013】

【課題を解決するための手段】上記目的を達成するため、第1の発明は、入力話者の音声から、標準話者の音声スペクトルに対して上記入力話者の音声スペクトルの周波数軸を伸縮する際の周波数伸縮関数を話者特徴として抽出する話者特徴抽出装置において、所定の音声単位毎に、上記標準話者の音響モデルに対して、上記入力話者の音声サンプルの尤度あるいは音響モデルを上記入力話者の音声サンプルに話者適応させた話者適応音響モデルの尤度を最大にするという基準に従って、上記周波数伸縮関数を最尤推定し、この推定された上記周波数伸縮関数の集合の頻度分布を求める頻度計測手段と、上記頻度分布に基づいて、最大頻度を有する周波数伸縮関数を話者特徴として抽出するモード抽出手段を備えたことを特徴としている。

【0014】上記構成によれば、頻度計測手段によって、所定の音声単位毎に、上記周波数伸縮関数の最尤推定が行われ、この推定された周波数伸縮関数の集合の頻度分布が求められる。そして、モード抽出手段によって、最大頻度を有する周波数伸縮関数が話者特徴として抽出される。したがって、上記周波数伸縮関数の頻度分布に複数のピークが存在しても、正確な周波数伸縮関数が安定して抽出される。

【0015】また、上記第1の発明の話者特徴抽出装置は、入力話者の音声サンプルから音響モデルを用いたビタピアルゴリズムによって音素境界情報を推定する音素境界情報推定手段を備えて、上記頻度計測手段を、上記音素境界情報に基づいて、有音無音の別および調音点の位置に従って、上記最尤推定を行う際に用いる上記入力話者の音声サンプルの音声区間を限定するように成すこ

とが望ましい。

【0016】上記構成によれば、音素境界情報推定手段で推定された音素境界情報に基づいて、上記頻度計測手段によって、有音無音の別および調音点の位置に従って上記入力話者の音声サンプルの音声区間が限定される。したがって、声道長の差の影響を受け難い音素や無音部を上記最尤推定時における周波数伸縮の対象外にして、声道長の差の影響を受け難い音素や無音部まで変形されることが防止可能になる。こうして、少ない発声データから、発声データの内容に依存せずに安定して話者特徴が抽出される。

【0017】また、上記第1の発明の話者特徴抽出装置は、音響モデルの各状態が、有音無音の別および調音点の位置に従って予め設定された上記最尤推定を行う音声区間に属しているか否かを判別し、属している状態に関して、上記音響モデルを入力音声サンプルに話者適応させて話者適応音響モデルを作成する話者適応モデル作成手段を備えて、上記頻度計測手段を、上記標準話者音響モデルに対して上記話者適応音響モデルを用いて上記最尤推定を行うように成すことが望ましい。

【0018】上記構成によれば、話者適応モデル作成手段によって、有音無音の別および調音点の位置に従って設定された音声区間に属する音響モデルの状態に関して、入力音声サンプルに話者適応された話者適応音響モデルが作成される。そして、上記頻度計測手段によって、上記標準話者音響モデルに対して上記話者適応音響モデルを用いて上記最尤推定が行われる。したがって、声道長の差の影響を受け難い音素や無音部を上記話者適応音響モデル作成の対象外にして、声道長の差の影響を受け難い音素や無音部までが変形されるのを防止することが可能になる。こうして、少ない発声データから、発声データの内容に依存せずに安定して話者特徴が抽出される。

【0019】また、上記第1の発明の話者特徴抽出装置は、上記モード抽出手段を、上記最大頻度を有する周波数伸縮関数が複数存在する場合には、上記頻度分布を混合ガウス分布で表現した場合における当該複数の周波数伸縮関数が属している分布の分散が大きい方の周波数伸縮関数をもって話者特徴するように成すことが望ましい。

【0020】上記構成によれば、上記モード抽出手段によって、上記最大頻度を有する周波数伸縮関数が複数存在する場合には、属している分布の分散が大きい方の周波数伸縮関数が抽出される。こうして、より多くの周波数伸縮関数の集団を代表する周波数伸縮関数が話者特徴として抽出される。

【0021】また、上記第1の発明の話者特徴抽出装置は、上記モード抽出手段を、上記標準話者の特徴を表す周波数伸縮関数に近い方の周波数伸縮関数をもって話者特徴とする機能を組み合わせて、上記話者特徴を抽出するよ

うに成すことが望ましい。

【0022】上記構成によれば、上記最大傾度を有する複数の周波数伸縮関数が属している分布の分散の大きさが同程度である場合には、より標準話者の周波数伸縮関数に近い方の周波数伸縮関数が、話者特徴として抽出される。

【0023】また、第2の発明は、標準話者の音声スペクトルに対して入力話者の音声スペクトルの周波数軸を伸縮する際の周波数伸縮関数を用いて入力話者の音声スペクトルの周波数軸を伸縮することによって上記入力話者の音声を正規化する正規化手段を有する音声認識装置において、上記正規化手段は、上記第1の発明の話者特徴抽出装置と、上記話者特徴抽出装置によって抽出された周波数伸縮関数を用いて、上記入力話者の音声スペクトルの周波数軸を伸縮する周波数ワープ手段で構成されていることを特徴としている。

【００２４】上記構成によれば、周波数ワープ手段によって、上記第１の発明の話者特徴抽出装置で抽出された話者の特徴をよりの確に表す正確な周波数伸縮関数を用いて話者正規化が行われる。したがって、発声データの内容に依存せずに安定して話者正規化が行われて、高い精度で認識結果が得られる。

【0025】また、第3の発明は、入力話者の音声スペクトルに対して標準話者の音声スペクトルの周波数軸を伸縮する際の周波数伸縮関数を用いて音声のスペクトルの周波数軸を伸縮することによって音響モデルを入力話者に話者適応させる話者適応手段を有する音声認識装置において、上記話者適応手段は、上記第1の発明の話者特徴抽出装置と、上記話者特徴抽出装置によって抽出された周波数伸縮関数の逆関数を用いて、上記音響モデルの周波数軸を伸縮する周波数ワープ手段で構成されていることを特徴としている。

【００２６】上記構成によれば、周波数ワープ手段によって、上記第１の発明の話者特徴抽出装置で抽出された話者の特徴をよりの確に表す正確な周波数伸縮関数の逆関数を用いて、話者適応が行われる。したがって、発声データの内容に依存せずに安定して話者適応が行われて、高い精度で認識結果が得られる。

【0027】また、第4の発明は、入力話者の音声スペクトルに対して標準話者の音声スペクトルの周波数軸を伸縮する際の周波数伸縮関数を用いて音声のスペクトルの周波数軸を伸縮することによって、標準話者の音声素片を接続して成る合成音声の声質を発話者の声質に変換する声質変換手段を有する音声合成装置において、上記声質変換手段は、上記第1の発明の話者特徴抽出装置と、上記話者特徴抽出装置によって抽出された周波数伸縮関数の逆関数を用いて、上記音声素片の周波数軸を伸縮する周波数ワープ手段で構成されていることを特徴としている。

【００２８】上記構成によれば、周波数ワープ手段によ

として、上記第1の発明の話者特徴抽出装置で抽出された話者の特徴をよりの確に表す正確な周波数伸縮関数の逆関数を用いて、声質変換が行われる。したがって、発声データの内容に依存せずに安定して声質変換が行われて、より入力話者の声質に近い合成音声を得られる。

【0029】また、第5の発明は、入力話者の音声から、標準話者の音声スペクトルに対して上記入力話者の音声スペクトルの周波数軸を伸縮する際の周波数伸縮関数を話者特徴として抽出する話者特徴抽出方法において、所定の音声単位毎に、上記標準話者の音響モデルに対して、上記入力話者の音声サンプルの尤度あるいは音響モデルを上記入力話者の音声サンプルに話者適応させた話者適応音響モデルの尤度を最大にするという基準に従って、上記周波数伸縮関数を最尤推定し、この推定された上記周波数伸縮関数の集合の頻度分布を求め、上記頻度分布に基づいて、最大頻度を有する周波数伸縮関数を話者特徴として抽出することを特徴としている。

【0030】上記構成によれば、所定の音声単位毎に最尤推定された周波数伸縮関数の集合の頻度分布に基づいて、最大頻度を有する周波数伸縮関数が話者特徴として抽出される。したがって、上記周波数伸縮関数の頻度分布に複数のピークが存在しても、正確な周波数伸縮関数が安定して抽出される。

【0031】また、第6の発明のプログラム記録媒体は、コンピュータを、上記第1の発明における頻度計測手段およびモード抽出手段として機能させる話者特徴抽出処理プログラムが記録されたことを特徴としている。

【００３２】上記構成によれば、上記第１の発明の場合と同様に、上記周波数伸縮関数の頻度分布に複数のピークが存在する場合でも、正確な周波数伸縮関数が安定して抽出される。

【 0 0 3 3 】

【発明の実施の形態】以下、この発明を図示の実施の形態により詳細に説明する。

＜第1実施の形態＞図1は、本実施の形態の音声認識装置におけるブロック図である。尚、この音声認識装置は、話者正規化方式を用いた音声認識装置であり、上記HMMに代表される音響モデルをベースとしている。

【0034】音声入力部1において、マイクから入力された音声はデジタル波形に変換されて音響分析部2に入力される。音響分析部2は、入力されたデジタル波形を短い時間間隔(フレーム)毎に周波数分析し、スペクトルを表す音響パラメータのベクトル系列に変換する。ここで、上記周波数分析としては、MFCC(メル周波数FFT(高速フーリエ変換)ケプストラム)やLPC(線形予測分析)メルケプストラム等のスペクトルを効率よく表現できる音響パラメータを抽出できる分析方法が用いられる。こうして得られた音響パラメータ系列は、話者正規化部3を構成する周波数ワープ部4に送出され

る。

【0035】上記話者正規化部3は、上記周波数ワーブ部4と周波数ワーブ関数推定部5とから概略構成される。そして、周波数ワーブ関数推定部5は、学習時には、音響分析部2からの音響パラメータ系列と単語列入力部6から入力された学習用単語の音素列とに基づいて、音素境界情報および周波数ワーピング関数を推定して周波数ワーブ部4に送出する。また、認識時には、音響分析部2からの音響パラメータ系列と上記学習時に推定された周波数ワーピング関数とに基づいて音素境界情報を推定し、この推定された音素境界情報を上記周波数ワーピング関数と共に周波数ワーブ部4に送出する。尚、周波数ワーブ関数推定部5の構成と動作については後に詳述する。

【0036】そうすると、上記周波数ワーブ部4は、上記周波数ワーピング関数および音素境界情報を用いて、入力音声の音響パラメータ系列を周波数ワーブ(話者正規化)し、周波数ワーブ後の音響パラメータ系列を尤度演算部7に送出するのである。そして、尤度演算部7では、周波数ワーブされた音響パラメータ系列に対して、不特定話者音響モデル格納部8に格納された不特定話者モデル(HMM)を作用させて、各音韻の状態毎に尤度を算出する。そして、得られた尤度系列を照合部9に送出する。

【0037】上記照合部9は、上記尤度演算部7からの尤度系列に対して、辞書格納部10に登録された総ての言語モデル(単語)との照合を行ない、各単語のスコアを算出する。そして、上位のスコアを呈する単語を認識候補(認識結果)として出力部11から出力するのである。

【0038】以下、上記周波数ワーブ関数推定部5の構成と動作について詳細に詳述する。図2に、上記周波数ワーブ関数推定部5における学習時に機能する部分の構成を示す。さらに、図3には、周波数ワーブ関数推定部5における認識時に機能する部分の構成を示す。先ず、図2に従って、学習時について説明する。

【0039】音素境界推定部15は、全話者音響モデル格納部12に格納された混合数1以上の全話者音響モデル(HMM)を用いて、ビタビアルゴリズムによって音素境界情報を求める。その際に、教師あり学習時には、音素境界推定部15には、単語列入力部6からの音素列と音響分析部2からの音響パラメータ系列(学習データ)とが入力される。そうすると、音素境界推定部15は、入力音響パラメータ系列に入力音素列を適用させて、上記全話者音響モデルを用いたビタビアルゴリズムによって音素境界情報を求める。これに対して、教師なし学習時には、音素境界推定部15には、音響分析部2からの音響パラメータ系列(学習データ)のみが入力される。そうすると、音素境界推定部15は、入力音響パラメータ系列に言語モデル格納部13に格納された弱い文法の言語モデルを適用させて、全話者音響モデルを用いたビタビアルゴリズムによって音素境界情報を求めるのである。

そして、こうして得られた音素境界情報は頻度計測部16および周波数ワーブ部4に送出される。

【0040】尚、上記「弱い文法」とは、対象言語の音素または音節の接続に関する制約条件のみを表現するネットワーク(有限状態オートマトン)のことである。例えば、日本語の場合には、/k/と/i/とは接続するが、/s/と/k/は接続しないというような制約条件である。また、上記音素境界情報とはこの音素境界情報によって分離される音素のラベル情報をも含む概念であり、上記ビタビアルゴリズムによって求まる。

【0041】上記頻度計測部16は、入力された音素境界情報に従って、全話者音響モデル格納部17に格納された混合数1の全話者モデル(HMM)を用いて、後に詳述する方法によって、音響モデルの状態や入力サンプル等の音声単位毎に周波数ワーピング関数 f の係数 α を最尤推定する。さらに、 α 軸に関する頻度の分布を表す分布関数 $H(\alpha)$ を求める。そして、得られた分布関数 $H(\alpha)$ をモード抽出部18に送出する。

【0042】上記モード抽出部18は、後述するようにして、上記分布関数 $H(\alpha)$ の中から最大頻度を与える最適係数 α (つまり周波数ワーピング関数 f)を推定する。そして、推定された周波数ワーピング関数 f を関数格納部19に格納すると共に、周波数ワーブ部4に送出するのである。

【0043】次に、図3に従って、認識時について説明する。尚、図3における全話者音響モデル格納部12、言語モデル格納部13、音素境界推定部15および関数格納部19は、図2において学習時に使用される全話者音響モデル格納部12、言語モデル格納部13、音素境界推定部15および関数格納部19と同じものである。

【0044】事前ワーブ部20は、上記学習時に推定されて関数格納部19に格納された周波数ワーピング関数 f を用いて、認識対象の音響パラメータ系列を周波数ワーブする。以下、この場合の周波数ワーブを、後に周波数ワーブ部4によって行われる周波数ワーブに対して「事前ワーブ」と言うことにする。こうして、事前ワーブが行われた音響パラメータ系列が音素境界推定部15に送出される。

【0045】そうすると、上記音素境界推定部15は、事前ワーブが行われた音響パラメータ系列に弱い文法の言語モデルを適用させて、全話者音響モデルを用いたビタビアルゴリズムによって音素境界情報を求めるのである。その場合、認識対象の音響パラメータ系列は、学習時に抽出された話者特徴としての周波数ワーピング関数 f を用いて事前ワーブされている。したがって、より話者の声道長に即した音素境界情報を求めることができるのである。そして、得られた音素境界情報が、関数格納部19に格納されている周波数ワーピング関数 f と共に周波数ワーブ部4に送出される。

【0046】そうすると、上記周波数ワーブ部4におい

ては、上記推定された周波数ワーピング関数 f によって、上記学習時には、入力された学習用の音響パラメータ系列が周波数ワープされる。一方、上記認識時には、入力された認識用の学習音響パラメータ系列が周波数ワープされるのである。

【0047】すなわち、本実施の形態においては、上記学習時における周波数ワープ関数推定部 5 による係数 α の最尤推定を、全音声単位に関して行うのではなく個々の音声単位毎に行い、その最大頻度を呈する係数 α を推定することによって、正確な周波数ワーピング関数 f を安定して推定するのである。また、上記学習時に頻度計測部 16 と周波数ワープ部 4 とによって入力音響パラメータ系列に周波数ワーピング関数 f を適用する場合、および、認識時に周波数ワープ部 4 によって入力音響パラメータ系列に周波数ワーピング関数 f を適用する場合には、後に詳述するように、表 1 の分類表に従って、上記音素境界情報に基づいて、周波数ワープ(正規化)の対象とする音素区間を限定するのである。こうすることによって高精度認識を行う音声認識装置を構築することができるのである。

【0048】ところで、上記周波数ワープ関数推定部 5 における上記周波数ワーピング関数 f の推定方法には、以下に述べる二通りの推定方法がある。

(A) 標準話者の音響モデルを入力音声データに話者適応させた適応モデルを用いる。

(B) 入力音声データを直接用いる。

そして、この二通りの推定方法を、入力音声データの量や質に応じて使い分けるのである。ここで、音声データの質とは尤度の上昇具合であり、周波数ワープ関数推定部 5 は、上記二通りの推定方法による尤度の上昇具合を見計らって、上昇の大きい推定方法を採用するのである。長いエンロール期間を許容できる音声認識装置の場合には、このような推定処理も可能となる。尚、長いエンロール期間を許容できない場合には、予め何れかの推定方法に固定しておけばよい。

【0049】上記推定方法(A)は、入力音声データが少ない場合に有効である。また、推定方法(B)は、入力音声データが多い場合に有効であり、入力音声データから直接求めるために、精密な推定が可能となる。但し、入力音声データが少ない場合には、当該推定をエンロールモードで行う際に、入力音声データに無い音素環境における上記係数の推定や平滑化が問題になる。

【0050】また、上記推定方法(A)、(B)の各々に関して、使用する音響モデルは、全話者モデルの場合と、話者クラスタ別に作成された混合数が 1 の音響モデルの場合との二通りがある。音声認識装置の記憶容量が少ない場合には前者を採用する。一方、記憶容量が多い場合は音響モデル群を各話者クラスタ別に格納できるので後者を採用する。後者の場合には、入力音声データに基づいて最適な話者クラスタを選択し、この選択話者クラス

タに属する音響モデルを使用することになる。すなわち、図 1 に示す音声認識装置は、全話者モデルを用いた推定方法(B)によって上記係数の推定を行うのである。

【0051】ここで、上記話者クラスタ別に作成された音響モデルとは、ある基準で全学習話者をクラスタリングしておき、複数の話者クラスタ毎に学習によって作成された音響モデルのことである。ここでは、上記クラスタリングの基準として、各話者の声道情報を用いる。尚、周波数ワープ関数推定部 5 が使用する際には、適切な話者クラスタの音響モデルを選択して用いることになる。

【0052】次に、上記周波数ワープ関数推定部 5 が学習時に使用する全話者モデルと、尤度演算部 7 が上記尤度演算時に使用する不特定話者モデルの作成方法について説明する。上記全話者モデルは、総ての学習話者の音声データを用いて学習した音響モデルである。通常、周波数ワーピング関数 f の最尤推定に使用する場合には混合数を 1 に設定する。これに対して、不特定話者モデルは、通常学習話者の音声データをそのまま用いて学習した音響モデルである。しかしながら、本実施の形態のように話者正規化を行う音声認識装置においては、尤度演算部 7 に入力される音響パラメータ系列は、周波数ワープ部 4 によって既に正規化されている。したがって、不特定話者モデルも、学習話者の音声データを以下に述べる正規化と同様の手順で正規化した正規化学習データを用いた学習によって作成するのである。その場合、不特定話者に対応させるために、通常では、混合数は 1 以上に設定される。

【0053】次に、上記周波数ワープ関数推定部 5 によって行われる周波数ワーピング関数 f の推定について説明する。まず、周波数ワーピング関数 f の定義について説明する。周波数ワーピング関数 f (周波数伸縮関数または単に伸縮関数と言う場合もある)の周波数軸は声道の長さを直接反映しているので声道長伸縮関数とも言う。周波数ワーピング関数 f は、推定の容易さを考慮して、通常はできるだけ少ないパラメータ数で表現される。本実施の形態においては、周波数ワーピング関数 f を、以下のようなパラメータが 1 個からなる区間線形関数であると定義する。

【0054】周波数ワーピング関数 $f()$:

・ $x \leq \min(\omega/\alpha, \omega)$ では、 $f(x) = \alpha x$

($\omega \approx 4 \text{ kHz}$) ($0.88 < \alpha < 1.13$)

・ $\min(\omega/\alpha, \omega) < x$ では、

$\alpha > 1$ のとき $f(x) \rightarrow (\omega/\alpha, \omega)$ と $(fs/2, fs/2)$ とを結ぶ直線

$\alpha \leq 1$ のとき $f(x) \rightarrow (\omega, \alpha\omega)$ と $(fs/2, fs/2)$ とを結ぶ直線

ここで、 α : 周波数ワーピング関数 $f()$ の係数

fs : サンプリング周波数

尚、上記サンプリング周波数 fs は、本実施の形態におい

ては8kHz以上を仮定している。すなわち、 $f_s = 12\text{kHz}$ の場合には、 $(f_s/2, f_s/2)$ は(6kHz, 6kHz)となるのである。また、係数 α の定義域「 $0.88 < \alpha < 1.13$ 」は飽くまでも一例であり、子供まで含めると「 $0.7 < \alpha < 1.13$ 」となる。 $\alpha > 1$ である場合における上述のような折れ線で表される周波数ワーピング関数 $f(x)$ を図4に示す。すなわち、周波数ワーピング関数 $f()$ の推定とは係数 α を推定することである。

【0055】また、複数のパラメータを有する周波数ワーピング関数 $f(x)$ の場合でも、以下に述べる1個のパラメータ α を有する周波数ワーピング関数 $f(x)$ の場合と同様に、パラメータ空間の総ての座標における尤度を算出して頻度を計測することによって、音響パラメータ系列に適合した係数を推定することができる。

【0056】上記周波数ワープ関数推定部5における上記周波数ワーピング関数 f の推定方法が上記推定方法(A)である場合には、標準話者の音響モデルを入力音声データに話者適応させた適応モデルを用いて、以下の手順によって2つの音響モデルの状態間の尤度を求めて推定するのである。

【0057】尚、その場合における上記適応モデルは、*

$$\hat{\alpha}_i = \arg \max_{\alpha} r_i(\hat{\mu}_i^f) \quad (i \in \Omega) \quad \dots(1)$$

ここで、 Ω ：正規化対象の出力確率密度関数集合のインデックス

$r_i()$ ：標準モデルの第 i 番目の出力確率密度関数

$f()$ ： α を係数とする周波数ワーピング関数

μ_i^f ：入力モデルにおける第 i 番目の出力確率密度関数 $b_i()$ の平均値ベクトル μ_i を $f()$ で周波数ワープしたベクトル

尚、上記正規化対象の出力確率密度関数集合 Ω は、後述する正規化対象の音素区間に属する音素に関する音響モデルの出力確率密度関数の集合である。

【0060】上記音響分析部2による音響分析で得られる音響パラメータや上記音響モデルの出力確率密度関数の引数は、通常MFCCやLPCケプストラムである。これらの音響パラメータの各次元はケプストラムと呼ばれる物理量であって、周波数ではない。そこで、上記周波数ワープ処理を行なう際には、学習データである音響パラメータからスペクトルへの変換 C^{-1} (ケプストラムの場合は逆cos変換)を行なって周波数次元に変換する。そして、周波数ワープ処理終了後は、逆変換 C (ケプストラムの場合はcos変換)を行なって元の音響パラメータ次元に戻すのである。すなわち、 $\mu_i^f = C(f(C^{-1}(\mu_i)))$ となる。ここで、 C^{-1} 、 C は、音響パラメータからスペクトルへの変換とその逆変換である。

【0061】次に、上記式(1)によって求められた $\{\hat{\alpha}_i\} (i \in \Omega)$ に関して、 α 軸に関する頻度の分布を求め、この頻度分布を表す関数を $h(\alpha)$ とおく。そして、上述のごとく正規化対象の出力確率密度関数集合 Ω の状態 i

*例えば、音響モデルの各状態が予め設定された正規化対象の音素区間に属しているか否かを判別し、属している状態に関して、上記音響モデルを入力音声データに話者適応させる話者適応モデル作成手段によって作成すればよい。

【0058】ここで、標準モデル(全話者音響モデル格納部17に格納された全話者モデルに相当)と入力モデル(上記適応モデルに相当)との2つの音響モデルの対応する状態間の尤度を、標準モデルの出力確率密度関数 $r_i()$ に、入力モデルの出力確率密度関数 $b_i()$ の平均値ベクトル μ_i を上記周波数ワーピング関数 $f()$ で周波数ワープして得られたベクトル μ_i^f を代入したときの値と定義する。上記各出力確率密度関数は多次元ガウス分布であって、平均値ベクトルと分散ベクトルから成っている。

【0059】そして、上記正規化対象の出力確率密度関数集合 Ω における第 i 番目の状態間の尤度 $r_i(\mu_i^f)$ に基づいて、状態 i における周波数ワーピング関数 $f()$ の最適係数 $\hat{\alpha}_i$ は、式(1)に示すように尤度 $r_i(\mu_i^f)$ を最大にする係数として推定されるのである。

毎に最尤推定して得られた係数 $\hat{\alpha}_i$ のうちモード(並数、最頻値)を与える係数 $\hat{\alpha}$ を、上記周波数ワーピング関数 $f()$ の最適係数として式(2)によって推定するのである。

$$\hat{\alpha} = \arg \max_{\alpha} h(\alpha) \quad \dots(2)$$

【0062】図5に、上記係数 α の頻度分布を表す分布関数 $h(\alpha)$ の一例を示す。このような多峰性を有する場合、すなわち複数のピークが存在する場合には、従来法によれば2つのピークの間にも最適係数 $\hat{\alpha}$ が求まる。これに対し、本実施の形態によれば、頻度が高い方のピークを呈するに係数 α が最適係数 $\hat{\alpha}$ として求まるのである。尚、図中、棒グラフは係数 α を0.01きざみで観測した場合の頻度分布 $h(\alpha)$ であり、破線はその包絡線である。サンプル数が少ない場合はこの包絡線を分布関数 $h(\alpha)$ として差し支えない。

【0063】ここで、図6に例示するように、上記分布関数 $h(\alpha)$ の最大値を与える α が複数個存在する場合がある。このような場合における最適係数 $\hat{\alpha}$ の推定は、「係数 α が属する分布の分散が大きいこと」および「標準話者に近いこと」の2つの基準を組み合わせで行う。図6においては、モード α_1 に属する分布は、分散が小さく、モード α_2 に比べて標準値の1.0から離れており、係数 α の推定誤りによるゴミと考えられる。

【0064】今、2つの値 α_1, α_2 で最大値 $h(\alpha)$ が与えられたとする。つまり、 $h(\alpha_1) = h(\alpha_2) = h(\hat{\alpha})$ となる場合である。係数 α が属する分布(すなわち、各

α_k を平均値とするガウス分布)の分散は、分布関数 $h(\alpha)$ を混合ガウス分布(α_k, σ_k^2)で表現することによって与えられる。混合ガウス分布の推定には、HMMの学習方法であるBaum-Welchアルゴリズムが用いられる。ここで、混合数は最大値を与える個数であり、図6の例の場合は「2」である。尚、 α_k は平均値であり、 σ_k^2 は分散であり、図6の場合は $k=1, 2$ である。

【0065】こうした場合、 α_1 を平均値とする分布の α_2 を平均値とする分布に対する分散の小ささの度合いは、例えば式(3)で表される。

$$s(1, 2) = \sigma_2^2 - \sigma_1^2 \quad \dots (3)$$

また、標準話者の係数 α は $\alpha=1$ であるから、 α_k の標 *

$$g(k_1, k_2) = \lambda * s(k_1, k_2) + (1 - \lambda) * d_{k_1} \quad \dots (5)$$

ここで、 λ はシミュレーション実験に基づいて与えられる重み係数であり、 $[0, 1]$ の間、例えば0.7等に設定される。

【0067】このように、本実施の形態においては、従来のように、全状態に関して係数 α を最尤推定するのではなく、個々の状態 i 毎に最尤推定して最大頻度を呈する係数 α を求めるのである。こうすることによって、各状態 i 毎の係数 α_i の集合における頻度分布に複数のピークが存在する場合でも、正確な周波数ワーピング関数 $f()$ を安定して推定できるのである。また、その際における各状態 i 毎の係数 α の最尤推定を、上記正規化対象の音素区間に属する音素に関してのみ行うことによって、少ない音声データによって、精度良く周波数ワーピング関数 $f()$ を推定できるのである。

【0068】一方、上記周波数ワーピング関数推定部5における上記周波数ワーピング関数 f の推定方法が上記推定方法(B)である場合には、音響分析部2からの入力音響パラメータ系列を直接用いて、以下の手順によって周波数ワーピング関数 $f()$ の最適係数 α を推定する。尚、上述したごとく、図1に示す音声認識装置における周波数ワーピング関数推定部5には、上記推定方法(B)が適用されている。したがって、以下の推定手順を行うことにな ※

$$\hat{\alpha}_j = \arg \max_{\alpha} P(X_j^f | W_j) \quad (j \in \Psi) \quad \dots (6)$$

上記(1)～(4)の処理を総ての入力音声サンプル $\{X_j\}$ ($j \in \Psi$) に対して実行して、各サンプル X_j 毎の最適係数に α_j を求める。以上の手順(1)から手順(4)までの処理は、頻度計測部16によって行われる。

【0073】(5) 上記求められた α_j を係数とする周波数ワーピング関数 $f_j()$ を用いて、ビタビアルゴリズムによって、各サンプル毎に音素境界情報が求められる。そして、全入力音声サンプル $\{X_j\}$ ($j \in \Psi$) のうち、各サンプル毎の音素境界情報に基づく正規化対象と ★

$$\hat{\alpha}_j = \arg \max_{\alpha} P(\bar{X}_j^f | W_j) \quad (j \in \Psi) \quad \dots (7)$$

【0075】(7) 全サンプルの正規化対象音素区間における $\{\alpha_j\}$ に関して頻度分布を求め、頻度分布を表す

*準話者への近さは例えば式(4)で表される。

$$d_k = |\alpha_k - 1| \quad \dots (4)$$

そして、通常は、式(3)における $s(1, 2)$ の値に基づいて最適係数 α^* を選択する。そして、 $s(1, 2)$ の値が非常に小さく、両者の分散が同程度と見なされる場合には、式(4)における d_k の値に基づいて最適係数 α^* を選択するのである。

【0066】または、以下のようにして最適係数 α^* を選択してもよい。すなわち、 α_{k1} の α_{k2} に対するスコア $g(k_1, k_2)$ を式(5)で定義する。そして、このスコアの値を α_1 と α_2 について求め、小さい方を最適係数 α^* として選択するのである。

※。ここで、入力音声サンプル X_j のインデックス j の全集合を Ψ とおく。

【0069】(1) α に初期値を代入する。ここで、 $X_j = \{x_j^-(t)\}$ ($t=1, 2, \dots, T_j$) であり、「 $x_j^-(t)$ 」は時刻(フレーム) t における音響パラメータベクトル、 T_j は音響パラメータ系列 X_j における最終時刻(最終フレーム)である。

【0070】(2) 上記入力音響パラメータ系列 X_j に α を係数とする周波数ワーピング関数 $f()$ を作用させることによって周波数ワープを行う。そして、ビタビアルゴリズムを用いて、上記周波数ワープが行われた入力音響パラメータ系列 X_j^f の標準モデル(全話者音響モデル格納部17に格納された全話者モデルに相当)に対する累積尤度 $P(X_j^f | W_j)$ を求める。ここで、 W_j は入力音響パラメータ系列 X_j^f の音素列である。

【0071】(3) 係数 α を、定義域 $0.88 < \alpha < 1.13$ 内において、例えば0.02きざみで移動させながら、上記(2)の処理を繰り返して累積尤度 P を求める。

【0072】(4) 上記累積尤度 P を最大とする α_j を求め、これを α_j とおく。すなわち、 α_j は、式(6)によって表される。

★なる音素区間の音響パラメータ系列の集合を $\{X_j^-\}$ とおく。この手順(5)による音素境界情報算出処理は音素境界推定部15で行われる。

【0074】(6) 上記正規化対象となる音素区間の音響パラメータ系列の集合 $\{X_j^-\}$ に関して、上記音素境界情報に基づく正規化対象の音素区間毎に、ビタビアルゴリズムによって累積尤度 $P(X_j^f | W_j)$ を求める。そして、式(7)によって、最適係数 α_j^* が求め直される。

分布関数を $H(\alpha)$ とおく。そして、上述した推定方法

(A)の場合と同様に、係数 α_j^* のうちモード(並数, 最頻

値)を与える係数 α^* を、上記周波数ワーピング関数 $f()$ の最適係数として式(8)によって推定するのである。

$$\hat{\alpha} = \arg \max_{\alpha} H(\alpha) \quad \dots(8)$$

上記手順(6)とこの手順(7)における頻度分布の算出とは頻度計測部16によって行われる。また、手順(7)における上記モードを与える係数 α^* の抽出はモード抽出部18によって行われる。

【0076】尚、上記頻度分布に同一値のピークが複数の存在する場合には、上記推定方法(A)において述べた方法と同様の方法によって最適係数を推定する。

【0077】このように、上記推定方法(B)の場合には、個々のサンプルj毎に係数 α_j を最尤推定し、最大頻度を呈する α を最適係数として求めることによって、各サンプルj毎の係数 α_j の集合における頻度分布に複数のピークが存在する場合でも正確な周波数ワーピング関数 $f()$ を安定して推定できるのである。また、その際における各サンプルj毎の係数 α_j の最尤推定を、上記正規化対象の音素区間に属する音素に関してのみ行うことによって、少ない音声データによって、精度良く周波数ワーピング関数 $f()$ を推定できるのである。

【0078】尚、上記周波数ワーピング関数 $f()$ 推定処理における音素境界推定部15と頻度計測部16とモード抽出部18との処理の区分は、上述に限定されるものではない。例えば、手順(2)および手順(6)におけるビ*

*タビ演算を、音素境界推定部15で行うようにしても差し支えない。

【0079】次に、上記学習時には頻度計測部16と周波数ワーブ部4とで、認識時には周波数ワーブ部4で周波数ワーブを行う際に、頻度計測部16および周波数ワーブ部4によって行われる上記音素境界情報に基づく対象音素区間の限定について説明する。

【0080】上述したように、学習時および認識時においては、周波数ワーブ関数推定部5の音素境界推定部15によって、入力話者の音響パラメータ系列あるいはこの入力音響パラメータ系列に基づく適応モデルに、発話内容の音素列や言語モデル格納部13に格納された弱い文法の言語モデルを適用させて、全話者音響モデル格納部12に格納された全話者モデルや話者クラスタにクラスタリングされた全話者モデルから選択されたものを用いたビタビアルゴリズムによって音素境界情報を求め、頻度計測部16(学習時)および周波数ワーブ部4(学習時、認識時)に送出するようにしている。

【0081】そうすると、上記頻度計測部16および周波数ワーブ部4は、上記周波数ワーブ関数推定部5からの音素境界情報に基づいて、入力音声データのうち周波数ワーブ処理の対象とする音素区間を制御するのである。本実施の形態においては、音素を表1に示す5種類に分類する。表1

分類	音素	説明
[a]	無音区間	環境騒音そのものであり、声道長とは全く関係ない。
[b]	調音点が歯茎より前に位置する子音	音源が声道出口付近に位置するために、声道全体による共鳴管が形成されず、声道長の影響を受け難い。
[c]	調音点が歯茎より後に位置する子音、半母音	母音と同様に声道長に直接影響される。しかしながら、子音や半母音区間は元来安定していないので声道長の推定には適さない。
[d]	「ウ」を除く母音	母音は全般的に声道長の影響を直接受ける。
[e]	母音「ウ」、撥音	日本語の「ウ」は発声の仕方によってフォルマンと周波数が大きく変動するので、声道長の影響よりも発声の仕方に大きく影響される。撥音も音素環境に大きく依存すると共に鼻音化の影響が大きい。

【0082】そして、この分類に基づいて、以下のような区別に従って、上記頻度計測部16は学習時の周波数ワーブを制御し、周波数ワーブ部4は学習時および認識時の正規化を制御するのである。

・学習時…分類[d]

・認識時…分類[c], 分類[d], 分類[e], (分類[b])

但し、認識時には、分類[b]を含めてもよい。発音の仕方によっては、音素「イ」も音素「ウ」と同様に狭母音なのでフォルマント周波数が大きく変動する場合がある。したがって分類[e]に音素「イ」を含め、分類[d]から音素「イ」を除いてもよい。

【0083】尚、上記周波数ワーブ部4による正規化処理対象の音素区間制御方法は、周波数ワーブ関数推定部

5の音素境界推定部15が用いる全話者モデルの規模に応じて二通りある。

・全話者音響モデル格納部12の容量に余裕がある場合には、全話者モデルの規模を非常に大きくできる場合には、分類[b]の調音点が歯茎より前に位置する子音を分離可能な音素境界情報を精度良く推定できるので、分類[c], 分類[d], 分類[e]のみを正規化対象区間とする。
 ・全話者モデルの規模をある程度大きくできる場合には、分類[b]を分離可能な音素境界情報を推定できないために上述のごとく分類[b]を入れて、分類[b], 分類[c], 分類[d], 分類[e]を正規化対象区間とする。つまり、無音区間のみを正規化対象の音素から外するのである。

【0084】上述したように、本実施の形態における音声認識装置は、高精度認識を行うために周波数ワープ部4において周波数ワープの対象とする音素区間を限定するようにしている。しかしながら、計算資源(処理能力)に余裕がないシステムに搭載する場合には、全話者モデルの規模を大きくできないため精度良く音素境界情報を推定することができない。そのような場合には、周波数ワープ部4を常に動作させて、全音素区間を対象に周波数ワープを行っても差し支えない。このように精度良く音素境界情報を推定できない場合でも、分類[d]の声道長の影響を直接受ける母音は推定できる。したがって、周波数ワープ関数推定部5の頻度計測部16によって推*

*定された周波数ワープ関数 $f()$ は、音素境界推定部15からの音素境界情報に基づいて声道長の影響を直接受ける分類[d]の母音のみから得られていることになる。したがって、周波数ワープ部4による周波数ワープの際に声道長の影響を受け難い音素区間と無音区間とが不必要に変形されることを防止するという効果は得ることができるのである。

【0085】最後に、上記周波数ワープ関数推定部5の音素境界推定部15が、学習時および認識時に用いる言語モデルについて説明する。表2に、各動作モード時における周波数ワープ関数推定部5が用いる言語モデルの切替状況を示す。表2

動作モード		言語モデル
学習	教師あり	単語列入力部6からの音素列
	教師なし	弱い文法の言語モデルまたは認識結果
認識	高精度	弱い文法の言語モデル
	通常	なし(不使用)

【0086】表2において、通常の認識処理時における言語モデル「なし」とは、上述のごとく全音素区間を正規化対象とするために正規化対象制御用の音素境界情報を推定する必要がなく、ビタビアルゴリズムを動作させないために言語モデルを使用しないという意味である。また、学習モードにおける「教師あり」とは、上述したように、音素境界情報の推定時にビタビアルゴリズムを行う際に発話内容の音素列を使用することであり、単語列入力部6から入力される音素列そのものが言語モデルとなる。これに対して、「教師なし」とは、発話内容の音素列を使用しないものであり、言語モデル格納部13に格納された弱い文法の言語モデルを使用するのである。

【0087】尚、上記弱い文法の言語モデルに代えて、認識結果を使用することも可能である。この場合、照合部9からの出力である認識結果を発話内容の音素列(言語モデル)として使用するのである。つまり、一度認識処理を行ってから再び学習モード時における周波数ワープ関数推定部5の処理動作に戻るのである。その場合には、図1に破線で示すように、出力部11からの認識単語列を一種の教師音素列として単語列入力部6に入力する。但し、発話内容に規制が無いので照合部9用の言語モデルを、音素境界推定部15でのビタビ演算に流用してよいかどうかという問題はあ

【0088】以上、上記HMMに代表される音響モデルを用いた音声認識装置について述べてきたが、音声波形または音声パラメータ系列を標準パターンとして登録しておくタイプの音声認識装置においても、入力音声データを直接用いる推定方法(B)の場合と同様の手法によってサンプル毎の頻度を観測することによって、本実施の形態における話者正規化方法を適用することができる。

$$\tilde{\alpha}_j = \arg \min_{\alpha} d(X_j^f | W_j) \quad \dots(9)$$

尚、その場合には、尤度の代わりにスペクトル間の距離尺度を用いることになる。処理手順は以下の通りである。

【0089】(1) α に初期値を代入する。入力音響パラメータ系列を $X_j = \{x_j(t)\}$ ($t=1, 2, \dots, T_j$)とおく。また、それに対応する標準パターンの音響パラメータ系列を $R_j = \{r_j(t)\}$ ($t=1, 2, \dots, T'_j$)とおく。なお、「 $x_j(t)$ 」、「 $r_j(t)$ 」は時刻(フレーム) t における音響パラメータベクトル、 j は各パターン(上記標準パターンに対応)のインデックス、「 T_j 」、「 T'_j 」は音響パラメータ系列 X_j 、 R_j における最終時刻(最終フレーム)である。

【0090】(2) 上記入力音響パラメータ系列 X_j に、 α を係数とする周波数ワープ関数 $f()$ を作用させて周波数ワープを行う。そして、周波数ワープの結果を X_j^f とおく。

【0091】(3) 上記周波数ワープ後の入力音響パラメータ系列 X_j^f と標準パターンの音響パラメータ系列 R_j との累積距離 $d(X_j^f, R_j)$ を、DPマッチングによって求める。尚、上記DPマッチングにおいては、距離尺度としてケプストラム距離等のスペクトル間距離を用いる。

【0092】(4) 係数 α を、定義域「 $0.88 < \alpha < 1.13$ 」内において、例えば0.02きざみで移動させながら、上記(2)と(3)との処理とを繰り返して累積尤度 d を求める。

【0093】(5) 上記累積尤度 d を最小とする α_j を求め、これを α_j^* とおく。すなわち、 α_j^* は、式(9)によって表される。

上記(1)～(5)の処理を総ての入力音声サンプル $[X_j]$ ($j \in \Psi$)に対して実行して、各サンプル X_j 毎の最適係数に α_j^- を求める。

【0094】(6) 上記求められた総ての α_j^- に関して頻度分布を求め、この頻度分布を表す分布関数を $H(\alpha)$ *

$$\hat{\alpha} = \arg \max_{\alpha} H(\alpha)$$

【0095】ここで、上記頻度分布に同一値のピークが複数の存在する場合には、上記推定方法(A)において述べた方法と同様の方法によって最適係数を推定する。

【0096】尚、サブワードHMMとは異なり、本例における上記標準パターンに音素情報は含まれていない。その場合には、音素等による正規化対象区間の制御は困難であるため導入はしない。その代わり、学習時には、一つのサンプルが単母音や特に表1における分類[d]に相当する母音で成る学習データを入力させるようにすることによって、正規化対象区間の制御を行えばよい。

【0097】上述したように、本実施の形態においては、上記音素境界推定部15、頻度計測部16およびモード抽出部18を有する周波数関数推定部5を備えている。そして、音素境界推定部15は、学習時には、音響分析部2からの入力音響パラメータ系列に、教師ありの場合には単語列入力部6からの音素列(言語モデル)を適用させる一方、教師なしの場合には言語モデル格納部13に格納された弱い文法の言語モデルを適用させて、全話者音響モデル格納部12に格納された全話者音響モデルを用いたビタビアルゴリズムによって音素境界情報を求める。

【0098】そうすると、上記頻度計測部16は、各サンプル j 毎に、上記周波数ワーピング関数 $f()$ の係数 α を定義域において所定値ずつ増加させながら入力音響パラメータ系列 X_j の周波数ワープを行う。そして、周波数ワープが行われた入力音響パラメータ系列 X_j^f のうち、上記音素境界情報に基づいて上記表1に従って上述のように設定された正規化対象区間の音響パラメータ系列のみに関して、全話者音響モデル格納部17に格納された全話者モデルに対する累積尤度 P を最大にする係数 α_j^- を最尤推定する。そして、各サンプル j 毎の $\{\alpha_j^-\}$ に関する頻度分布を表す分布関数 $H(\alpha)$ を求める。

【0099】さらに、上記モード抽出部18によって、係数 α_j^- のうち最頻値を与える係数 $\hat{\alpha}$ が周波数ワーピング関数 $f()$ の最適係数として推定され、この最適係数 $\hat{\alpha}$ を係数とする周波数ワーピング関数 $f()$ を関数格納部19に格納するのである。

【0100】これに対して、認識時には、上記音素境界推定部15によって、上記教師なし学習時と同様に弱い文法の言語モデルを適用させて、ビタビアルゴリズムによって音素境界情報を求めるのである。

【0101】こうして、上記周波数関数推定部5によって、推定された周波数ワーピング関数 $f()$ と音素境界情

*と置く。そして、上述した推定方法(A)、(B)の場合と同様に、係数 α_j^- のうち上記モードを与える係数 $\hat{\alpha}$ を上記周波数ワーピング関数 $f()$ の最適係数として式(10)によって推定するのである。

$$\dots(10)$$

報とが周波数ワープ部4に送出される。そして、周波数ワープ部4によって、上記音素境界情報に基づいて正規化対象となる音素区間が上記表1に従って学習時および認識時に応じて上述のように制御され、その制御結果に従って、当該認識対象の入力音響パラメータ系列が周波数ワープされるのである。

【0102】したがって、本実施の形態によれば、話者と標準話者との声道長の差を表わす声道長正規化係数 α を係数とする周波数ワーピング関数 $f()$ を用いて、頻度計測部16および周波数ワープ部4によって入力音響パラメータ系列を周波数ワープ(正規化)するに際して、周波数ワープの対象となる音素区間を制御することができる。その結果、声道長の差の影響を受け難い音素や無音部を正規化対象外とすることによって、声道長の差の影響を受け難い音素や無音部まで学習および正規化されてしまうことを防止できる。

【0103】さらに、上記学習時における周波数ワープ関数推定部5による係数 α の最尤推定を、全サンプル(推定方法(A)の場合には状態)に関して行うのではなく個々のサンプル(または状態)毎に行い、その頻度分布における最大頻度を呈する係数 α をもって周波数ワーピング関数 f の最適係数としている。したがって、上記頻度分布に複数のピークが存在する場合でも、正確な周波数ワーピング関数 f を安定して推定できるのである。

【0104】すなわち、本実施の形態においては、少ない発声データから安定して話者特徴を抽出し、その抽出結果を用いて精度よく話者正規化することによって、高い認識性能を得ることができるのである。

【0105】また、上記実施の形態においては、上記係数 α の分布関数 $h(\alpha)$ に最大値を与える α が α_1 と α_2 との2個存在する場合には、 α_1 を平均値とする分布の α_2 を平均値とする分布に対する分散の小ささの度合いを式(3)で求め、 α_k の標準話者への近さを式(4)で求める。そして、「係数 α が属する分布の分散が大きいこと」および「標準話者に近いこと」の2つの基準を組み合わせ、上記最適係数 α の推定を行うようにしている。したがって、同一最大ピーク値が複数存在するような分布関数 $h(\alpha)$ が得られた場合でも、安定して上記最適係数 α を推定することができるのである。

【0106】また、上記周波数ワープ部4による正規化対象となる音素区間の制御は、上記表1の音素分類に従って、学習時には分類[d]([f]を除く母音)を正規化対象音素区間とする。さらに、認識時には分類[c](調音

点が歯茎より後に位置する子音、半母音)、分類[d]、分類[e](母音「ウ」、撥音)、(分類[b](調音点が歯茎より前に位置する子音))を正規化対象音素区間とするようにしている。こうして、学習時および認識時における非正規化音素区間を、有音無音の別および調音点の位置に従って設定することによって、声道長の影響を受け難い音素区間と無音区間とが学習および正規化されることを、確実に防止することができるのである。

【0107】<第2実施の形態>図7は、本実施の形態の音声認識装置におけるブロック図である。尚、この音声認識装置は、話者適応方式を用いた音声認識装置である。音声入力部21、音響分析部22、単語列入力部26、尤度演算部27、照合部29、辞書格納部30および出力部31は、図1に示す上記第1実施の形態における音声入力部1、音響分析部2、単語列入力部6、尤度演算部7、照合部9、辞書格納部10および出力部11と同様である。また、周波数ワープ関数推定部24、全話者音響モデル格納部32、言語モデル格納部33および不特定話者音響モデル格納部34は、図1に示す周波数ワープ関数推定部5、全話者音響モデル格納部12、言語モデル格納部13および不特定話者音響モデル格納部8と同様である。尚、周波数ワープ関数推定部24、全話者音響モデル格納部32、言語モデル格納部33および不特定話者音響モデル格納部34は、周波数ワープ部25と共に、話者適応部23を構成している。

【0108】上記話者適応部23の周波数ワープ関数推定部24は、上記第1実施の形態の場合と同様にして、学習音響パラメータ系列に発話内容の音素列または弱い文法の言語モデルを適用して、全話者モデルを用いたビタビアルゴリズムを行って、音素境界情報および周波数ワープ関数 $f()$ を推定する。そうすると、周波数ワープ部25は、この推定された周波数ワープ関数 $f()$ の逆関数を用いて、不特定話者音響モデル格納部34に格納された不特定話者モデルを周波数ワープする。その場合、上記周波数ワープに際しては、上記音素境界情報に基づいて、上記表1における分類[b]、分類[c]、分類[d]、分類[e]に該当する音素の状態に対してのみ変換を行うことによって行う。そして、それ以外の状態は変換しないのである。但し、声道長の影響を受け難い分類[b]に該当する音素の状態は、変換しない場合もある。こうして周波数ワープされた不特定話者音響モデルを、話者適応モデル(HMM)として話者適応音響モデル格納部28に格納するのである。

【0109】こうして学習が終了すると、認識時には、上記尤度演算部27によって、音響分析部22からの入力音声の音響パラメータ系列に対して、話者適応音響モデル格納部28に格納された話者適応モデルを作用させて、上述した尤度演算処理を行なうのである。

【0110】このように、本実施の形態においては、学習時に、上記周波数ワープ関数推定部24によって、学

習音響パラメータ系列に基づいて上記音素境界情報および周波数ワープ関数 $f()$ を推定する。そして、周波数ワープ部25によって、上記推定された周波数ワープ関数 $f()$ の逆関数を用いて、分類[c]、分類[d]、分類[e]、(分類[b])に該当する音素の不特定話者モデルを周波数ワープすることによって、不特定話者モデルを話者適応させるようにしている。

【0111】したがって、本実施の形態によれば、上記不特定話者モデルを話者適応させる際における非正規化音素区間を、無音区間と長音点が歯茎より前に位置する子音とに設定することができる。その結果、声道長の影響を受け難い音素区間と無音区間とが不必要に変形されることを確実に防止することができるのである。

【0112】さらに、上記学習時における周波数ワープ関数 $f()$ の推定に際して係数 α の最尤推定を個々の状態やサンプル毎に行い、その最大頻度を呈する係数 α をもって周波数ワープ関数 f の最適係数としている。したがって、各状態やサンプル毎の係数 α の集合における頻度分布に複数のピークが存在する場合でも、正確な周波数ワープ関数 f を安定して推定することができる。

【0113】すなわち、本実施の形態によれば、少ない発声データから安定して話者特徴を抽出し、その抽出結果を用いて精度よく話者適応を行うことによって、高い認識性能を得ることができるのである。

【0114】尚、本実施の形態における上記話者適応音響モデル格納部28に格納する話者適応モデルの与え方には、上述の与え方の以外に、話者クラスタを用いる方法を採用してもよい。そして、この二通りの与え方を、音声認識装置の規模や入力音声データの量や質に応じて使い分けるのである。ここで、音声データの質とは尤度の上昇具合であり、話者適応部23は、上記二通りの与え方による尤度の上昇具合を見計らって、上昇の大きい推定方法を採用するのである。長いエンロール期間が許容できる音声認識装置の場合には、このような推定処理も可能となる。尚、上記話者クラスタを用いる方法においては、学習音声データに対する尤度が最大値になる話者クラスタの音響モデルを選択する。そして、この選択された音響モデルを話者適応モデルとして話者適応音響モデル格納部28に格納するのである。

【0115】また、上述した二つの与え方の何れかによって得られた話者適応モデルを初期モデルとして、上記MLLR方やVFS法等の既存の話者適応技術を用いて話者適応を行って新たに話者適応モデルを生成し、これを尤度演算部で用いるようにしても差し支えない。

【0116】<第3実施の形態>図8は、本実施の形態のテキスト音声合成装置におけるブロック図である。なお、このテキスト音声合成装置は、声質変換方式を用いたテキスト音声合成装置である。テキスト解析部41は、単語とそのアクセント型とが格納されたアクセント

辞書 42 を用い、入力テキストに対して形態素解析および係り受け解析を行って音素文字列とアクセント情報とを生成して韻律生成部 43 に送出する。韻律生成部 43 は、韻律制御テーブル 44 を参照して、継続時間長やピッチやパワーの韻律情報を生成して、音素文字列と共に音声素片選択部 45 に送出する。そうすると、音声素片選択部 45 は、音声素片辞書 46 から音素環境や韻律環境に最適な音声素片を選択し、音声素片情報を生成する。そして、この生成された音声素片情報を周波数ワー

部 48 に出力する一方、上記韻律情報を音声素片合成部 47 に出力する。

【0117】一方、周波数ワーブ関数推定部 49 は、声質変換のターゲット話者の入力音声波形を基に、第 1、第 2 実施の形態の場合と同様にして、上記音素境界情報および周波数ワーピング関数 $f()$ を推定する。そうすると、周波数ワーブ部 48 は、音声素片選択部 45 からの音声素片情報に含まれた音素境界情報に基づいて音質変換対象となる音素区間を上記表 1 に従って上述のように選択する。そして、その選択結果に従って、当該音質変換対象の音声素片情報である音響パラメータ系列を、上記推定された周波数ワーピング関数 $f()$ の逆関数を用いて周波数ワーブし、周波数ワーブ後の音声素片情報を音声素片合成部 47 に送出する。最後に、音声素片合成部 47 は、周波数ワーブ部 48 からの周波数ワーブ後の音声素片情報(音声素片の音響パラメータ系列)と音声素片選択部 45 からの韻律情報とを用いて、音声波形を生成しスピーカ 50 から音声出力するのである。

【0118】上述のように、本実施の形態においては、テキスト音声合成を行うに際して、上記周波数ワーブ関数推定部 49 によって、声質変換のターゲット話者における入力音声の音響パラメータ系列から上記音素境界情報および周波数ワーピング関数 $f()$ を推定する。そして、周波数ワーブ部 48 によって、上記音声素片情報に含まれた音素境界情報に基づいて音質変換対象となる音素区間を制御し、上記推定周波数ワーピング関数 $f()$ の逆関数を用いて、テキストに基づいて選択された音声素片の音質変換対象となる音響パラメータ系列を周波数ワーブすることによって、声質変換を行うようにしている。

【0119】したがって、本実施の形態によれば、テキストに基づいて選択された音声素片をターゲット話者の音質に変換する際における非声質変換音素区間を、無音区間と長音点が歯茎より前に位置する子音とに設定することができる。その結果、声道長の影響を受け難い音素区間と無音区間とが不必要に変形されることを確実に防止することができるのである。

【0120】さらに、上記学習時における周波数ワーピング関数 $f()$ の推定に際して係数 α の最尤推定を個々の状態やサンプル毎に行い、その最大頻度を呈する係数 α をもって周波数ワーピング関数 f の最適係数としてい

る。したがって、各状態やサンプル毎の係数 α の集合における頻度分布に複数のピークが存在する場合でも、正確な周波数ワーピング関数 f を安定して推定することができる。

【0121】すなわち、本実施の形態によれば、少ない発声データから安定して話者特徴を抽出し、その抽出結果を用いて精度よく声質変換を行うことによって正しく音質変換を行うことができるのである。

【0122】本実施の形態はスペクトル包絡の変換であり、声質の適応におおいに効果がある。しかしながら、話者間の声の特徴差は声質だけではなく韻律が大きく寄与する。したがって、本実施の形態に対して韻律の適応技術を併用しても構わない。

【0123】尚、上述した各実施の形態においては、上記周波数ワーブ部 4、25、48 において音響パラメータ系列を周波数ワーブする場合に、音声素片選択部 45 からの音声素片情報に含まれた音素境界情報に基づいて周波数ワーブの対象となる音素区間を制御するようにしている。しかしながら、この発明においては、必ずしもその必要はなく、総ての音素区間に対して周波数ワーブを行っても構わない。その場合であっても、周波数ワーブ関数推定部 5、24、49 によって推定された周波数ワーピング関数 $f()$ は、上記音素境界推定部 15 からの音素境界情報に基づいて声道長の影響を直接受ける分類[d]の母音のみから推定されている。したがって、周波数ワーブ部 4、25、48 による周波数ワーブの際に声道長の影響を受け難い音素区間と無音区間とが不必要に変形されることを防止するという効果は得ることができるのである。

【0124】また、上述した各実施の形態においては、上記周波数ワーピング関数 $f()$ で成る話者特徴を用いて話者正規化または話者適応を行う音声認識装置、および、上記周波数ワーピング関数 $f()$ で成る話者特徴を用いて声質変換を行う音声合成装置について説明している。しかしながら、この発明は、上記周波数ワーピング関数 $f()$ を話者特徴として抽出する話者特徴抽出装置にも適用されるものである。

【0125】ところで、その場合の話者特徴抽出装置における上記頻度計測手段およびモード抽出手段としての機能は、プログラム記録媒体に記録された話者特徴抽出処理プログラムによって実現される。上記プログラム記録媒体は、ROM(リード・オンリ・メモリ)でなるプログラムメディアである。あるいは、外部補助記憶装置に装着されて読み出されるプログラムメディアであってもよい。尚、何れの場合においても、上記プログラムメディアから話者特徴抽出処理プログラムを読み出すプログラム読み出し手段は、上記プログラムメディアに直接アクセスして読み出す構成を有していてもよいし、RAM(ランダム・アクセス・メモリ)に設けられたプログラム記憶エリア(図示せず)にダウンロードして、上記プログラ

ム記憶エリアにアクセスして読み出す構成を有しているもよい。尚、上記プログラムメディアからRAMの上記プログラム記憶エリアにダウンロードするためのダウンロードプログラムは、予め本体装置に格納されているものとする。

【0126】ここで、上記プログラムメディアとは、本体側と分離可能に構成され、磁気テープやカセットテープ等のテープ系、フロッピー（登録商標）ディスク、ハードディスク等の磁気ディスクやCD（コンパクトディスク）・ROM、MO（光磁気）ディスク、MD（ミニディスク）、DVD（デジタルビデオディスク）等の光ディスクのディスク系、IC（集積回路）カードや光カード等のカード系、マスクROM、EPROM（紫外線消去型ROM）、EEPROM（電氣的消去型ROM）、フラッシュROM等の半導体メモリ系を含めた、固定的にプログラムを担持する媒体である。

【0127】また、上記各実施の形態における音声認識装置、音声合成装置および話者特徴抽出装置は、モデムを備えてインターネットを含む通信ネットワークと接続可能な構成を有していれば、上記プログラムメディアは、通信ネットワークからのダウンロード等によって流動的にプログラムを担持する媒体であっても差し支えない。尚、その場合における上記通信ネットワークからダウンロードするためのダウンロードプログラムは、予め本体装置に格納されているものとする。または、別の記録媒体からインストールされるものとする。

【0128】尚、上記記録媒体に記録されるものはプログラムのみに限定されるものではなく、データも記録することが可能である。

【0129】

【発明の効果】以上より明らかなように、第1の発明の話者特徴抽出装置は、頻度計測手段によって、所定の音声単位毎に周波数伸縮関数を最尤推定して頻度分布を求め、モード抽出手段によって最大頻度を有する周波数伸縮関数を話者特徴として抽出するので、上記周波数伸縮関数の頻度分布に複数のピークが存在する場合でも、正確な周波数伸縮関数を安定して抽出することができる。したがって、この発明によれば、発声データの内容に依存せずに安定して話者特徴を抽出できる。

【0130】また、上記第1の発明の話者特徴抽出装置は、入力話者の音声サンプルから音素境界情報を推定する音素境界情報推定手段を備えて、上記頻度計測手段を、上記音素境界情報に基づいて、有音無音の別および調音点の位置に従って、上記最尤推定を行う際に用いる上記入力話者の音声サンプルの音声区間を限定するように成せば、声道長の差の影響を受け難い音素や無音部を上記最尤推定時における周波数軸伸縮の対象外にして、声道長の差の影響を受け難い音素や無音部まで変形されることを防止することが可能になる。したがって、少ない発声データから発声データの内容に依存せずに安定し

て話者特徴を抽出できる。

【0131】また、上記第1の発明の話者特徴抽出装置は、音響モデルの各状態が、有音無音の別及び調音点の位置に従って予め設定された上記最尤推定を行う音声区間に属しているか否かを判別し、属している状態に関して、上記音響モデルを入力音声サンプルに話者適応させて話者適応音響モデルを作成する話者適応モデル作成手段を備えて、上記頻度計測手段を、標準話者音響モデルに対して上記話者適応音響モデルを用いて上記最尤推定を行うように成せば、声道長の差の影響を受け難い音素や無音部を上記話者適応音響モデル作成の対象外にして、声道長の差の影響を受け難い音素や無音部まで変形されることを防止することが可能になる。したがって、少ない発声データから発声データの内容に依存せずに安定して話者特徴を抽出できる。

【0132】また、上記第1の発明の話者特徴抽出装置は、上記モード抽出手段を、上記最大頻度を有する周波数伸縮関数が複数存在する場合には、上記頻度分布を混合ガウス分布で表現した場合における当該複数の周波数伸縮関数が属している分布の分散が大きい方の周波数伸縮関数をもって話者特徴するように成せば、より多くの周波数伸縮関数の集団を代表する周波数伸縮関数を話者特徴として抽出することができる。

【0133】また、上記第1の発明の話者特徴抽出装置は、上記モード抽出手段を、上記標準話者の特徴を表す周波数伸縮関数に近い方の周波数伸縮関数をもって話者特徴とする機能を組み合わせて、上記話者特徴を抽出するように成せば、上記最大頻度を有する複数の周波数伸縮関数が属している分布の分散の大きさが同程度である場合でも、より適切な周波数伸縮関数を話者特徴として抽出することができる。

【0134】また、第2の発明の音声認識装置は、正規化手段を、上記第1の発明の話者特徴抽出装置と、上記話者特徴抽出装置によって抽出された周波数伸縮関数を用いて入力話者の音声スペクトルの周波数軸を伸縮する周波数ワープ手段で構成したので、話者の特徴をよりの確に表す正確な周波数伸縮関数を用いて話者正規化を行うことができる。したがって、発声データの内容に依存せずに安定して話者正規化を行って高い精度で認識結果を得ることができる。

【0135】また、第3の発明の音声認識装置は、話者適応手段を、上記第1の発明の話者特徴抽出装置と、上記話者特徴抽出装置で抽出された周波数伸縮関数の逆関数を用いて音響モデルの周波数軸を伸縮する周波数ワープ手段で構成したので、話者の特徴をよりの確に表す正確な周波数伸縮関数を用いて話者適応を行うことができる。したがって、発声データの内容に依存せずに安定して話者適応を行って高い精度で認識結果を得ることができる。

【0136】また、第4の発明の音声合成装置は、声質

変換手段を、上記第1の発明の話者特徴抽出装置と、上記話者特徴抽出装置によって抽出された周波数伸縮関数の逆関数を用いて標準話者の音声素片の周波数軸を伸縮する周波数ワープ手段で構成したので、話者の特徴をよりの確に表す正確な周波数伸縮関数を用いて声質変換を行うことができる。したがって、発声データの内容に依存せずに安定して声質変換を行って、より入力話者の声質に近い合成音声を得ることができる。

【0137】また、第5の発明の話者特徴抽出方法は、所定の音声単位毎に周波数伸縮関数を最尤推定し、この推定された上記周波数伸縮関数の集合の頻度分布を求め、最大頻度を有する周波数伸縮関数を話者特徴として抽出するので、上記周波数伸縮関数の頻度分布に複数のピークが存在する場合でも、正確な周波数伸縮関数を安定して抽出することができる。したがって、この発明によれば、発声データの内容に依存せずに安定して話者特徴を抽出できる。

【0138】また、第6の発明のプログラム記録媒体は、コンピュータを、上記第1の発明における頻度計測手段およびモード抽出手段として機能させる話者特徴抽出処理プログラムが記録されているので、上記第1の発明の場合と同様に、上記周波数伸縮関数の頻度分布に複数のピークが存在する場合でも、正確な周波数伸縮関数を安定して抽出することができる。したがって、発声データの内容に依存せずに安定して話者特徴を抽出できる。

【図面の簡単な説明】

【図1】 この発明の話者正規化方式を用いた音声認識装置におけるブロック図である。

【図2】 図1における周波数ワープ関数推定部の学習時に機能する部分の詳細なブロック図である。

【図3】 図1における周波数ワープ関数推定部の認識時に機能する部分の詳細なブロック図である。

【図4】 周波数ワーピング関数の一例を示す図であ

る。

【図5】 分布関数 $h(\alpha)$ の一例を示す図である。

【図6】 最大値を与える α が複数個存在する分布関数 $h(\alpha)$ を示す図である。

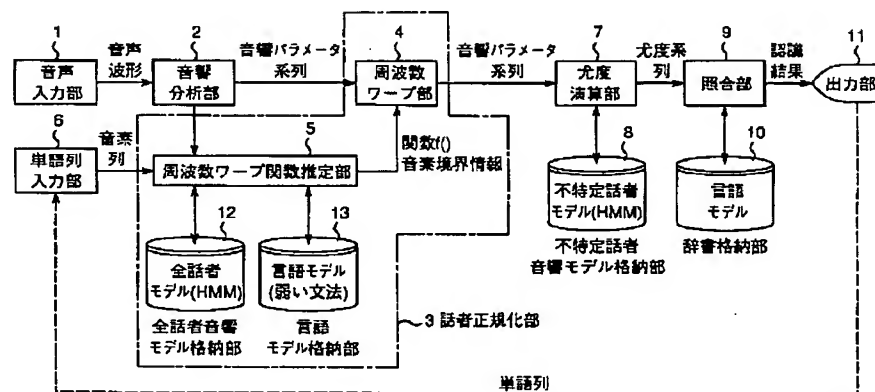
【図7】 図1とは異なる話者適応方式を用いた音声認識装置におけるブロック図である。

【図8】 この発明の音声合成装置のブロック図である。

【符号の説明】

- 1, 21…音声入力部、
- 2, 22…音響分析部、
- 3…話者正規化部、
- 4, 25, 48…周波数ワープ部、
- 5, 24, 49…周波数ワープ関数推定部、
- 6, 26…単語列入力部、
- 7, 27…尤度演算部、
- 8, 34…不特定話者音響モデル格納部、
- 9, 29…照合部、
- 10, 30…辞書格納部、
- 11, 31…出力部、
- 12, 17, 32…全話者音響モデル格納部、
- 13, 33…言語モデル格納部、
- 15…音素境界推定部、
- 16…頻度計測部、
- 18…モード抽出部、
- 19…関数格納部、
- 20…事前ワープ部、
- 23…話者適応部、
- 28…話者適応音響モデル格納部、
- 41…テキスト解析部、
- 43…韻律生成部、
- 45…音声素片選択部、
- 47…音声素片合成部、
- 50…スピーカ。

【図1】



[illegible]

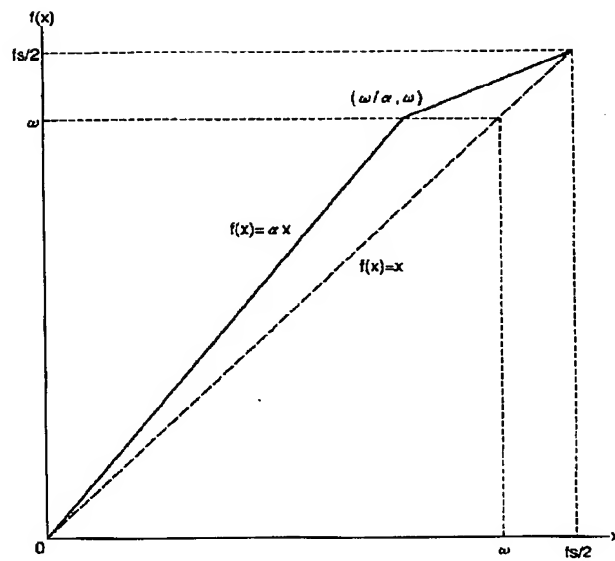
```

graph TD
    10[音楽パラメータ系列] --> 15[音楽境界推定部<br/>(ビタビアルゴリズム)]
    12[(全話者音声モデル<br/>多混合HMM)] --> 15
    13[(言語モデル<br/>弱い文法)] --> 15
    15 -- "関数f" --> 19[関数格納部]
    19 -- "関数f" --> 20[事前ワーブ部]
    20 --> 5[周波数ワーブ関数推定部]
    5 -- "音楽境界情報" --> 16[周波数ワーブ部4へ]
    19 -- "関数f" --> 16
  
```

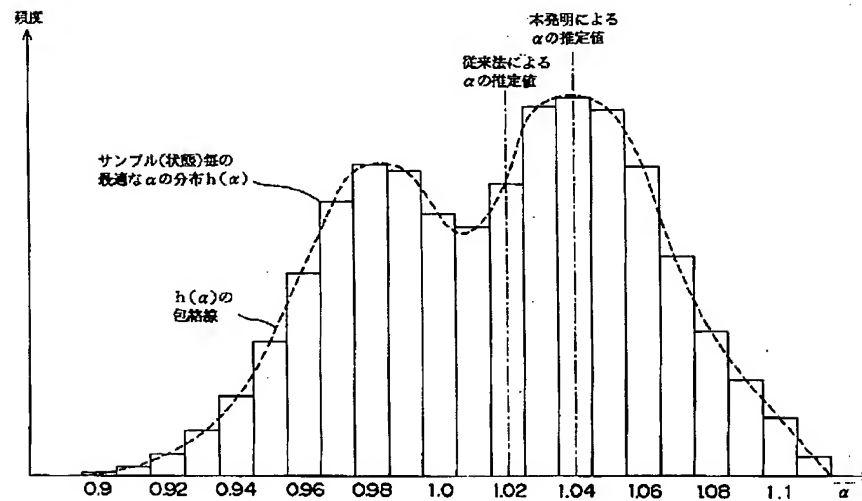
Figure 1 is a block diagram of a speech processing system. The system includes the following components and data flow:

- 21 音声入力部** (Voice Input Unit) receives input and outputs **22 音声波形** (Voice Waveform).
- 22 音響分析部** (Acoustic Analysis Unit) receives the waveform and outputs **音響パラメータ系列** (Acoustic Parameter Series).
- 27 尤度演算部** (Likelihood Calculation Unit) receives the parameter series and outputs **28 尤度系列** (Likelihood Series).
- 29 照合部** (Matching Unit) receives the likelihood series and outputs **31 認識結果** (Recognition Result).
- 31 出力部** (Output Unit) receives the recognition result.
- 28 単語列入力部** (Word Sequence Input Unit) outputs **28 音素列** (Phoneme Sequence).
- 24 関数f()音素境界情報** (Function f() Phoneme Boundary Information) receives the phoneme sequence and outputs **25 関波数ワープ部** (Function Waveform Warping Unit).
- 25 関波数ワープ部** outputs to **28 話者適応モデル(HMM)** (Speaker Adaptation Model (HMM)).
- 28 話者適応モデル(HMM)** is part of the **23 音響モデル格納部** (Acoustic Model Storage Unit).
- 28 話者適応モデル(HMM)** outputs to **29 照合部**.
- 30 言語モデル** (Language Model) is part of the **30 辞書格納部** (Dictionary Storage Unit) and outputs to **29 照合部**.
- 32 全話者モデル(HMM)** (All-Speaker Model (HMM)) is part of the **全話者音響モデル格納部** (All-Speaker Acoustic Model Storage Unit) and outputs to **24**.
- 33 言語モデル(語い文法)** (Language Model (Grammar)) is part of the **言語モデル格納部** (Language Model Storage Unit) and outputs to **24**.
- 34 不特定話者モデル(HMM)** (Unspecified Speaker Model (HMM)) is part of the **不特定話者音響モデル格納部** (Unspecified Speaker Acoustic Model Storage Unit) and outputs to **25**.

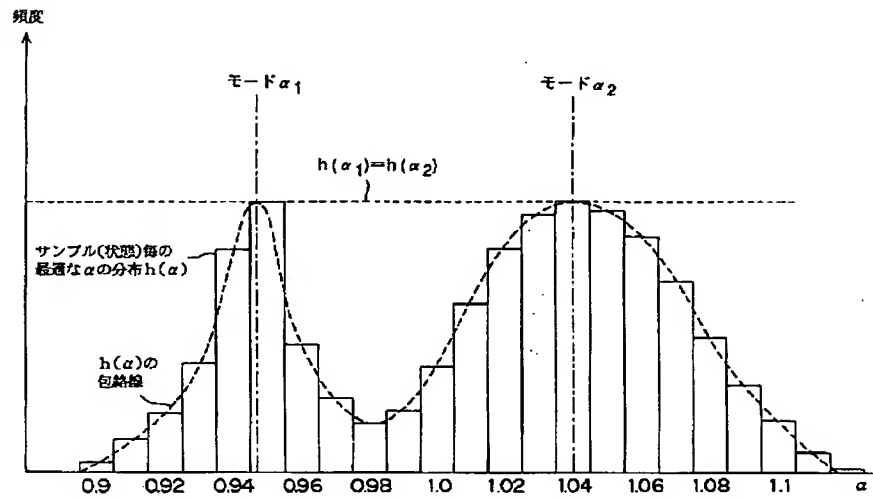
【図4】



【図5】



【図 6】



【図 8】

